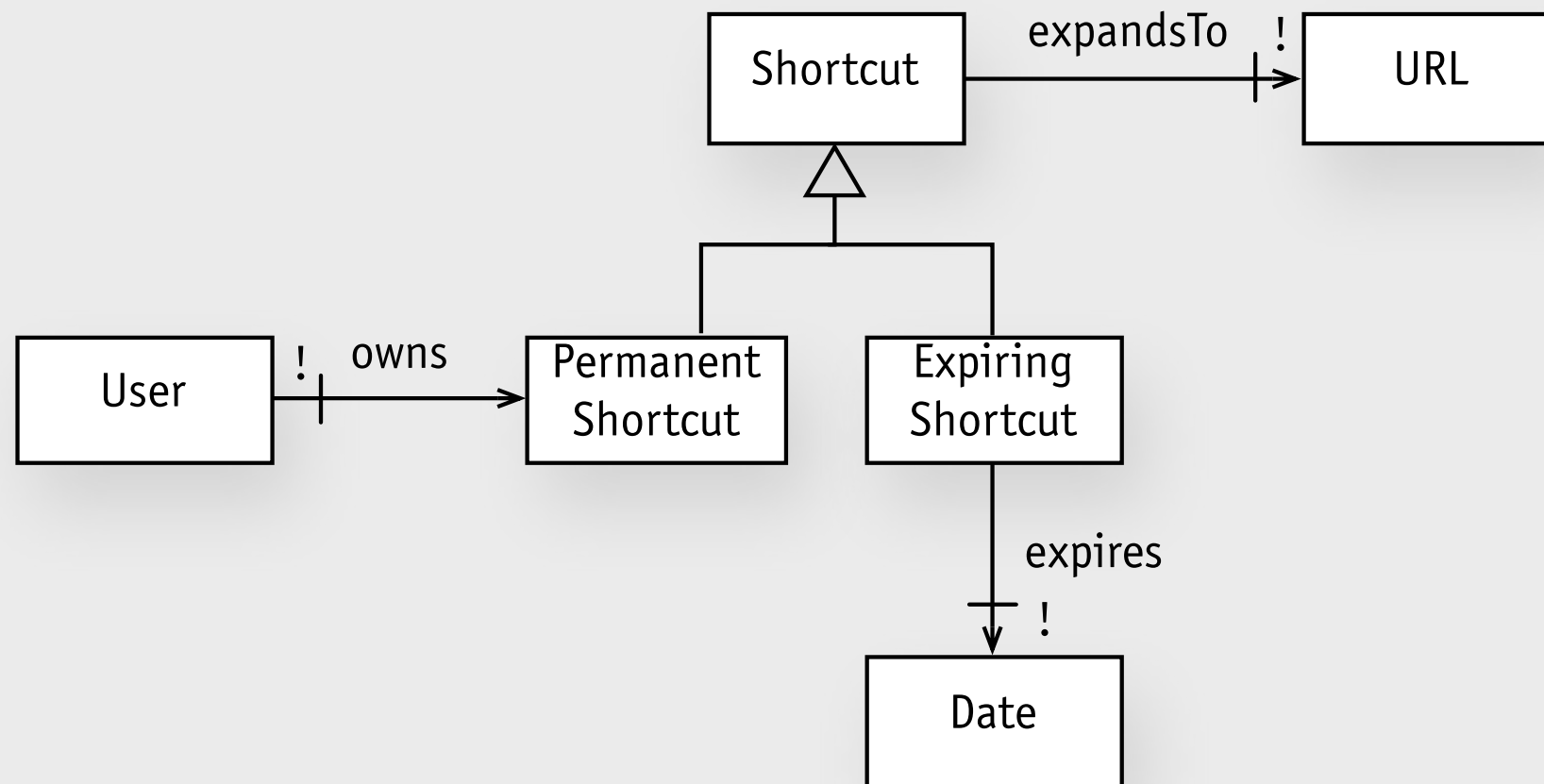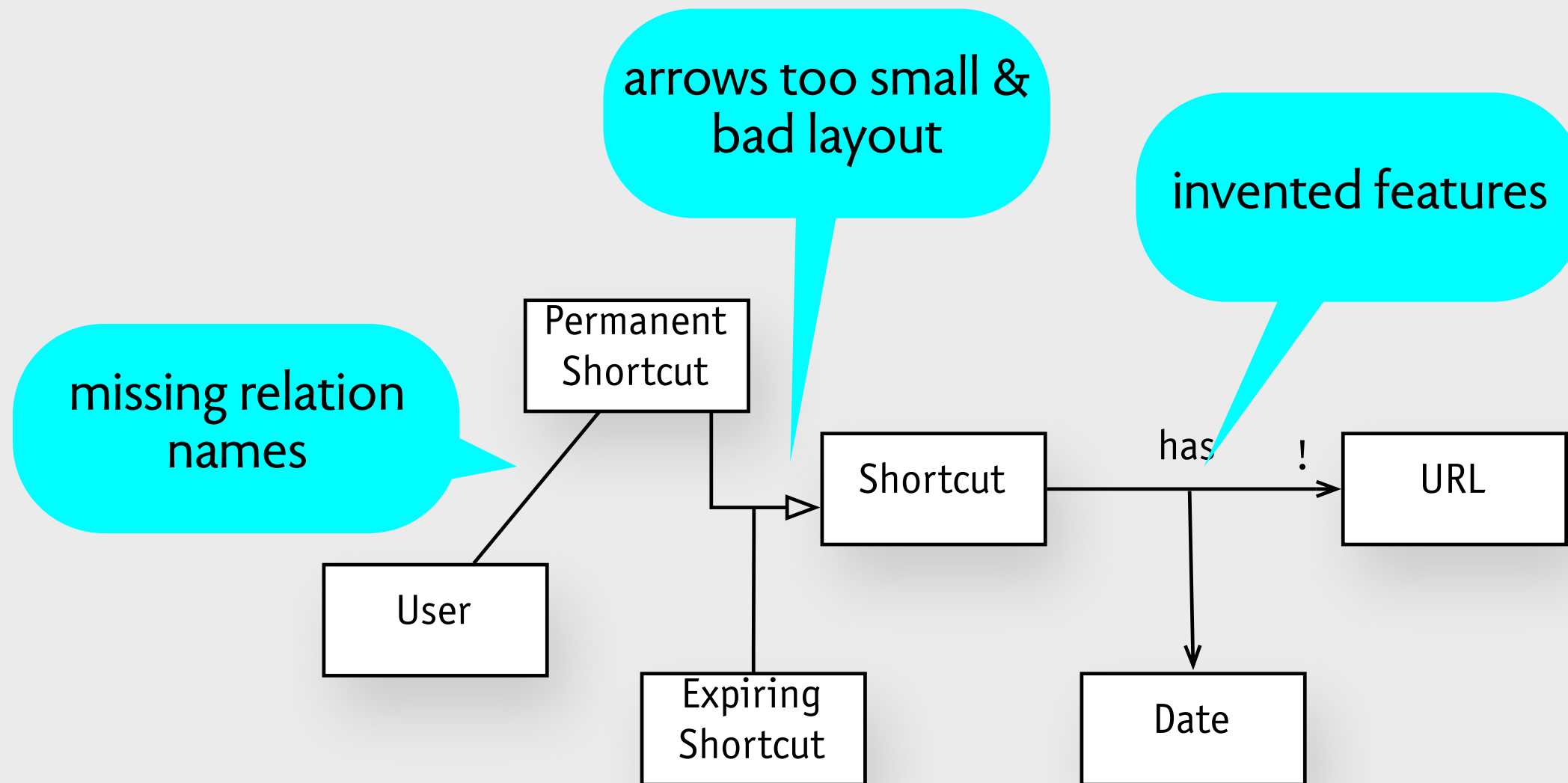# software studio

# data modeling errors

Daniel Jackson and Arvind Satyanarayan

classic errors
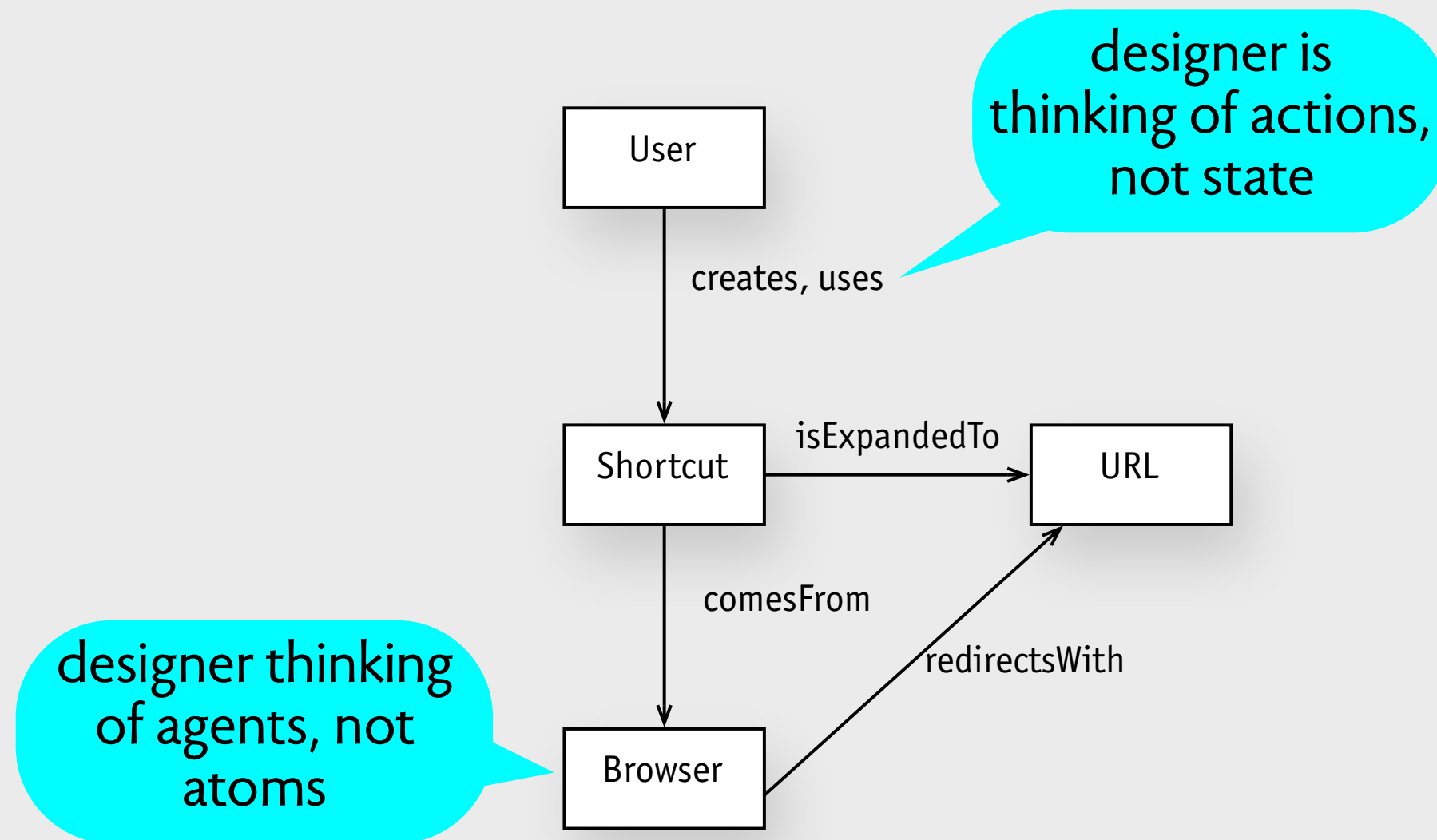
# suppose this is the data model

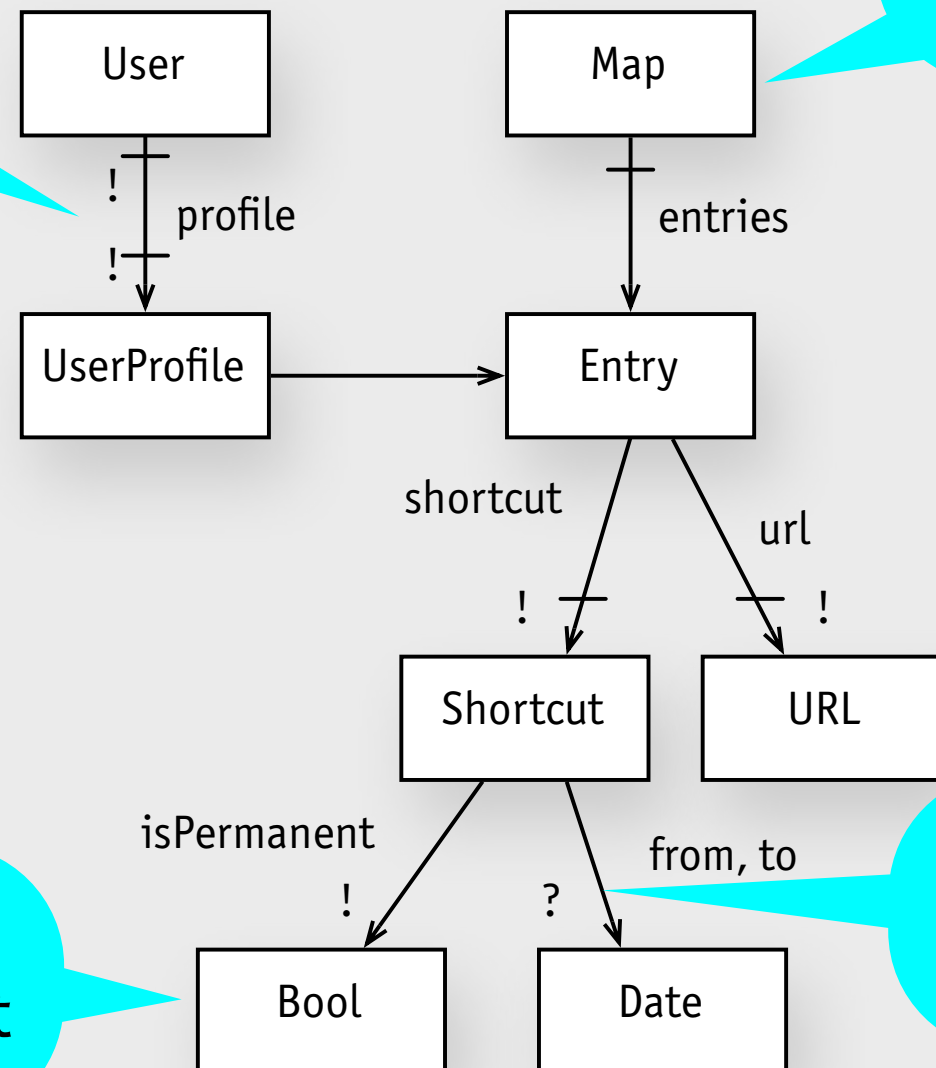# bad: not a data model syntactically

# bad: not a data model semantically

# bad: wrong semantics

# bad smells of data modeling

# don't be lazy with names



why?
it's never as obvious as it seems
choosing a good name helps designer get clarity

# beware primitive types

✗ [Address] —zip→ [String]

✓ [Address] —zip→ [Zipcode]

why?
type has syntactic or semantic properties
so may want to store in special way,
and/or use special validators

# don't mention low level ids

❌   | User |   --id-->   | UserId |

✔   | User |

why?
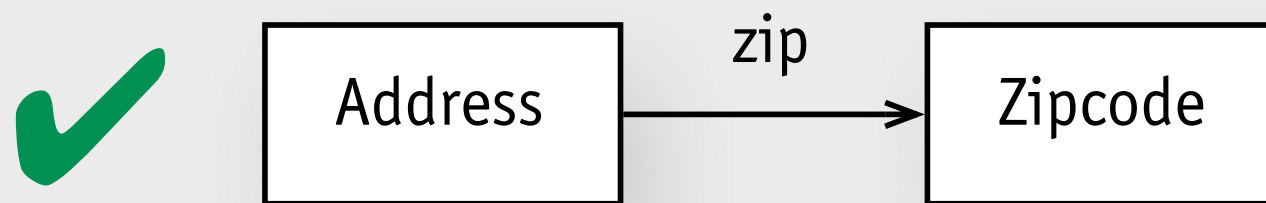every object has an implicit identity anyway
how it's represented is an implementation detail
but note: user-visible ids (such as usernames) *are* relevant

# don't split types

❌

✔

Movie —title→ MovieTitle
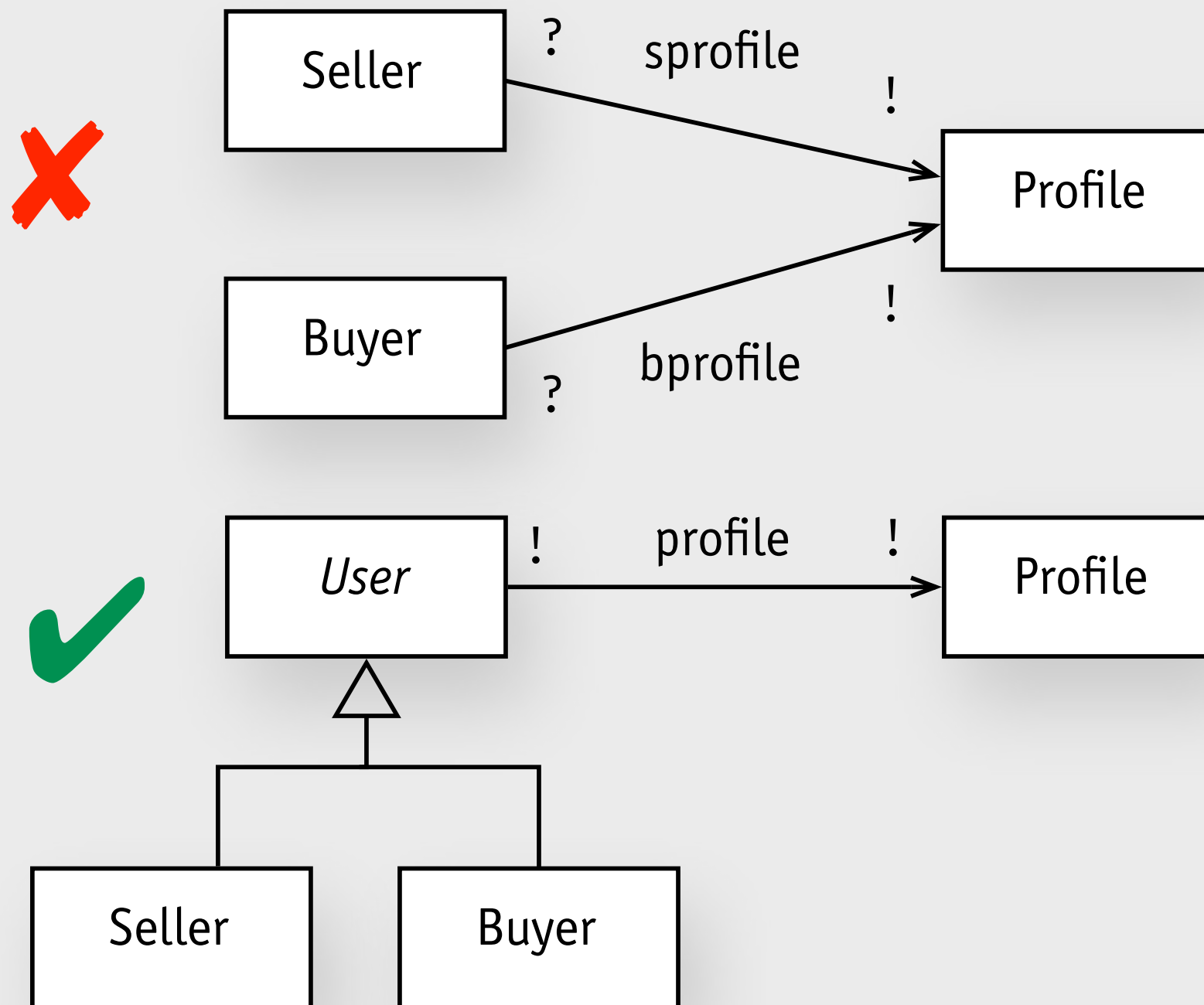
Movie —basedOn→ Book

Book —title→ BookTitle

Movie —mtitle→ Title

Movie —basedOn→ Book —btitle→ Title

Item —title→ Title

Item △ Movie, Book

why?
distinct types are disjoint, so couldn't ask whether movie and book have same title
so atoms to be compared must have the same type

# don't use set when order matters
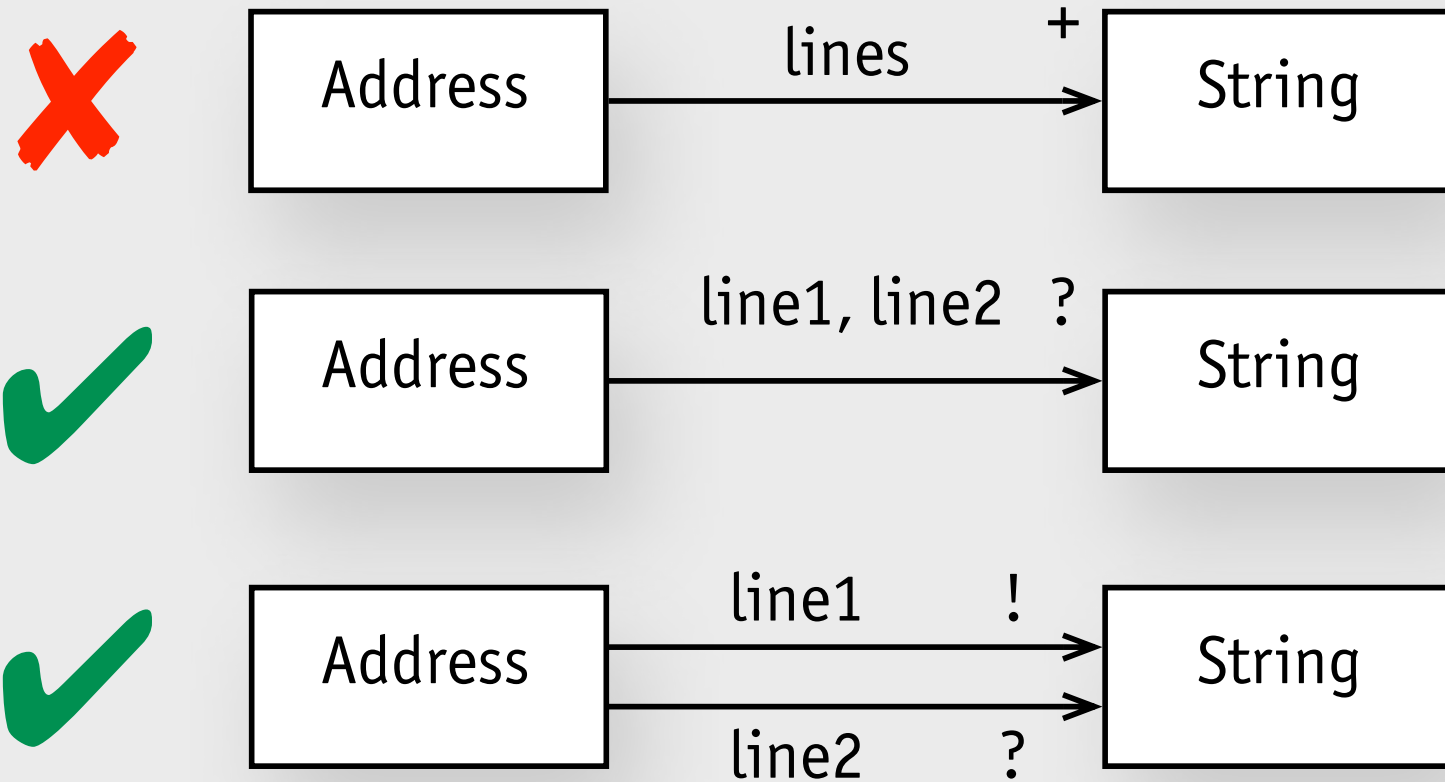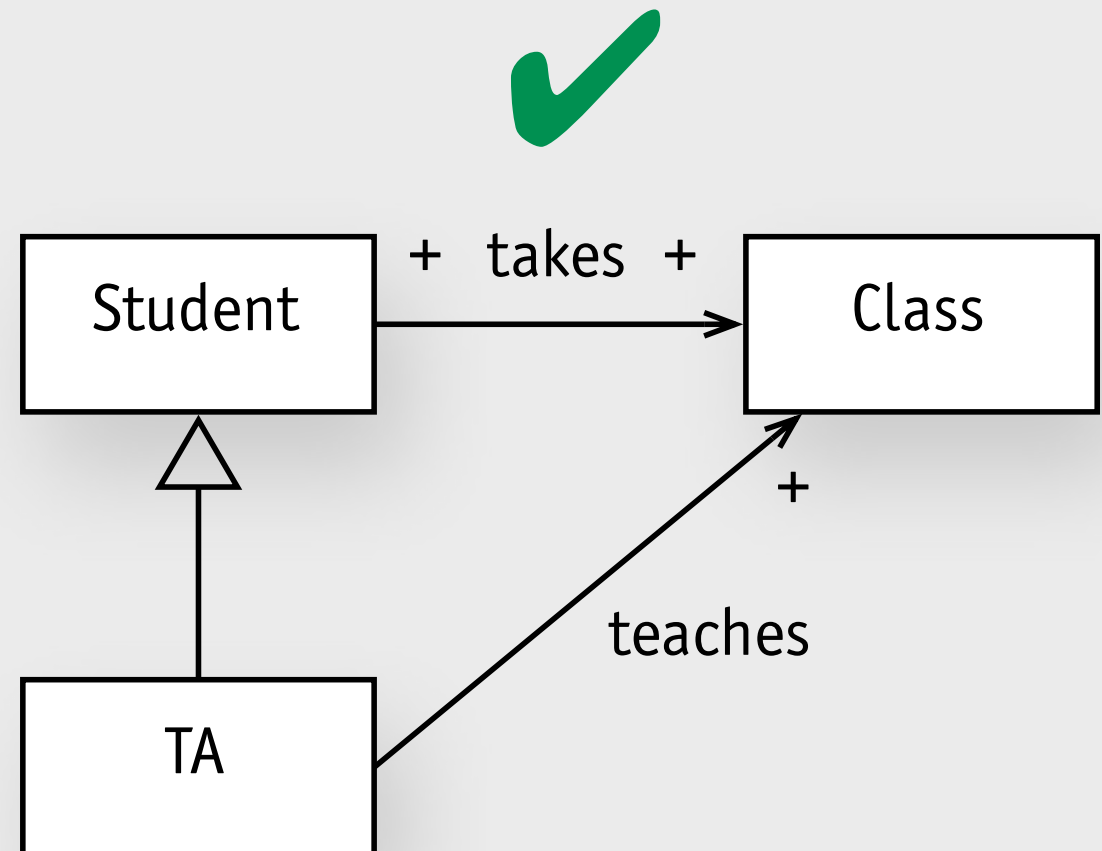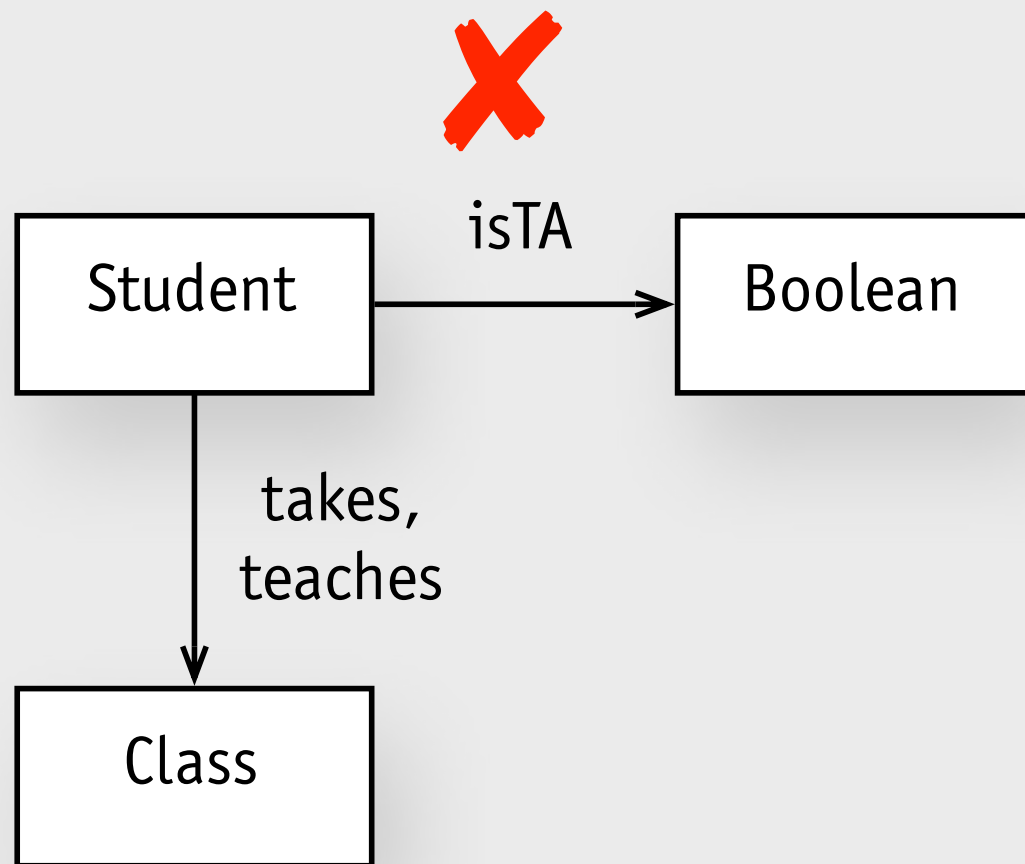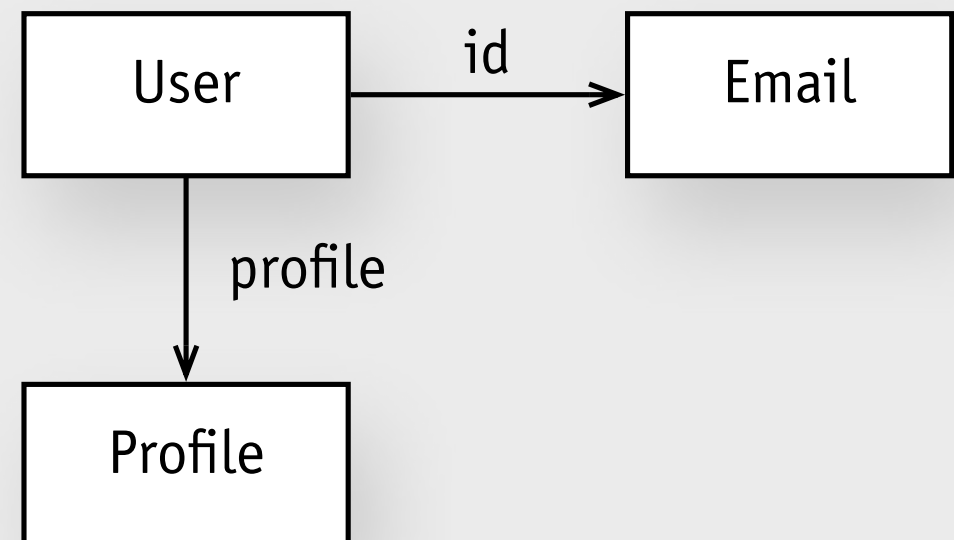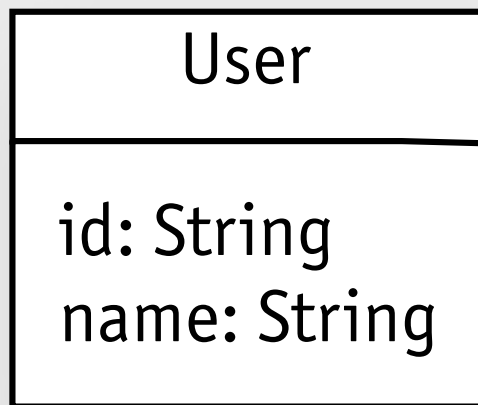


why?
tuples of a relation have no order
implementer can choose an unordered collection

# use subsets, not boolean flags

✘ ✔

Student —isTA→ Boolean

Student —+ takes +→ Class

Student ↑△ TA

TA —+ teaches→ Class

Student ↓takes, teaches→ Class

why?
flag is low level way to represent
obscures dynamic classification
prevents recording multiplicity graphically

# don't use attributes
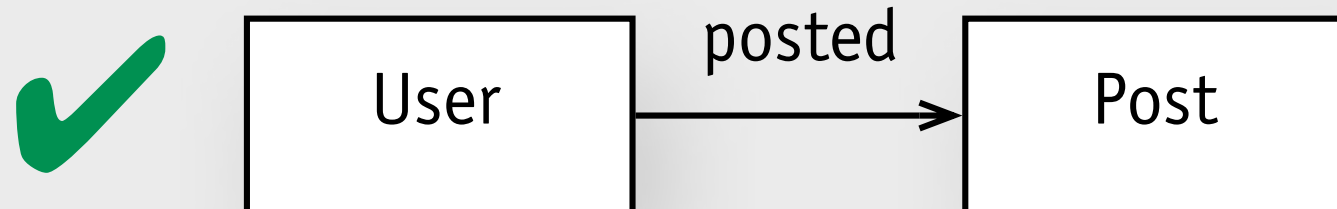
✗                                    ✔

| User |
|------|
| id: String<br>name: String |

```
┌──────┐   id   ┌───────┐
│ User │ ─────→ │ Email │
└──────┘        └───────┘
    │
 profile
    │
    ↓
┌─────────┐
│ Profile │
└─────────┘
```
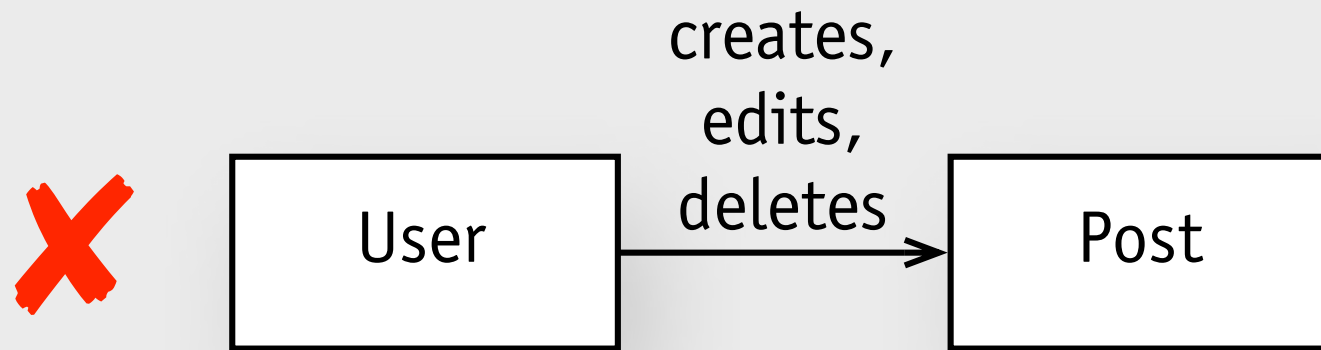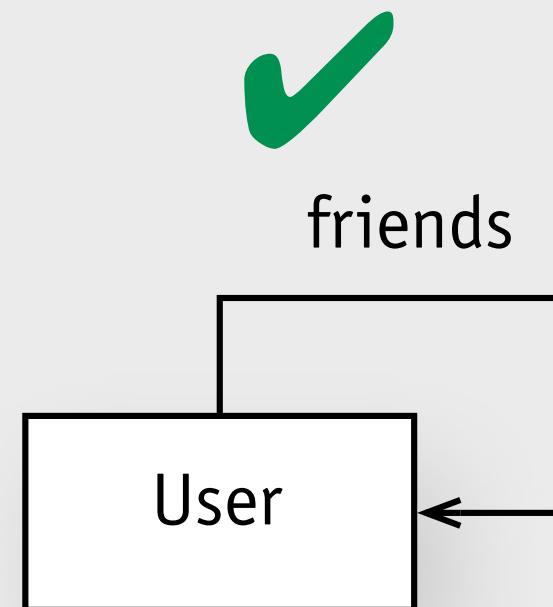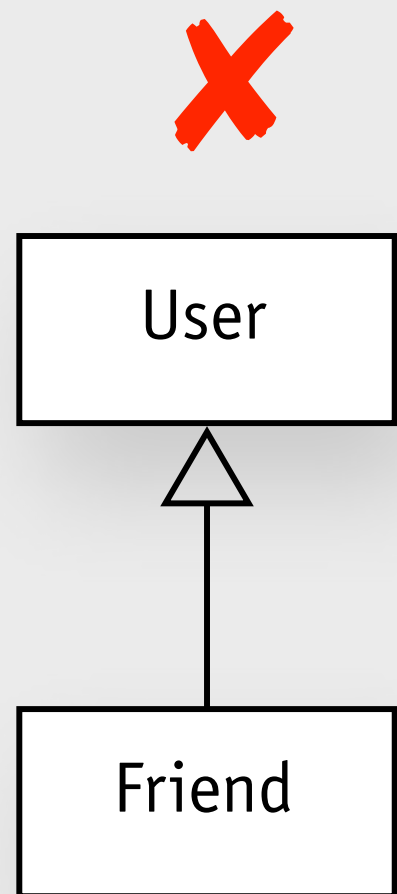
why?
attributes are an ill-defined idea, and just complicate the notation
better to factor out the relations that matter

# don't confuse state with actions

✗ User —creates, edits, deletes→ Post

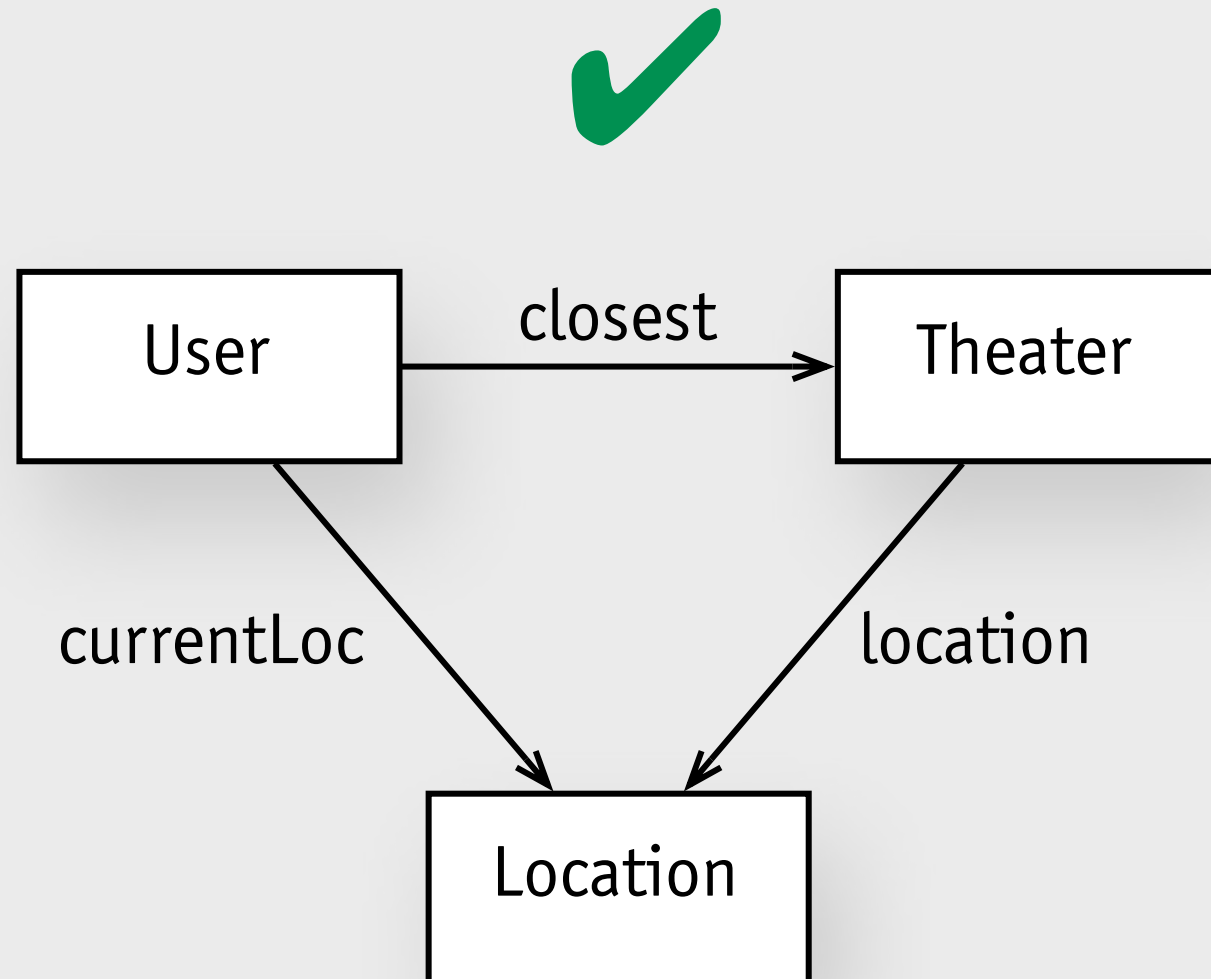✓ User —posted→ Post

why?
data model describes what is *remembered*
that is, what's stored in the state

# don't use subsets for relational state

✗

✔

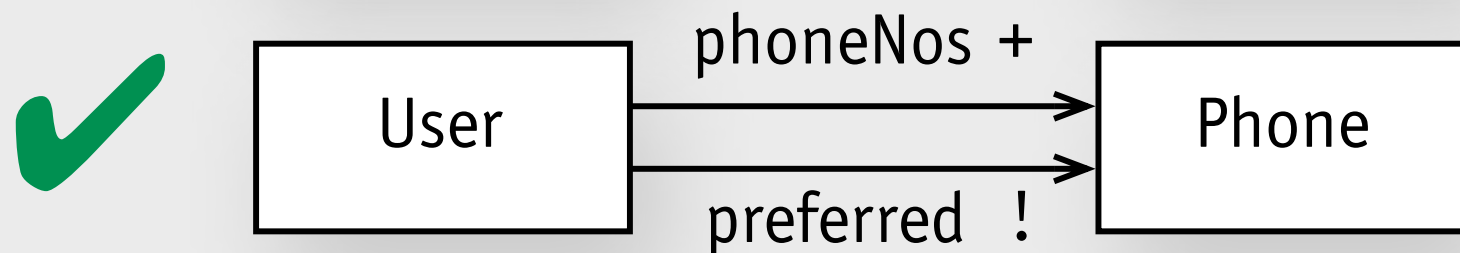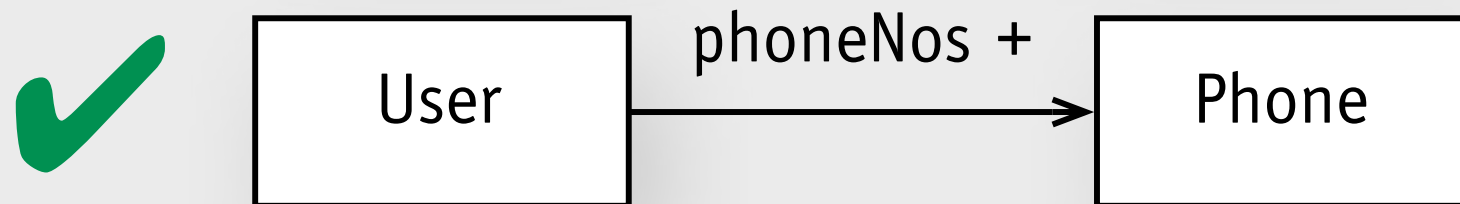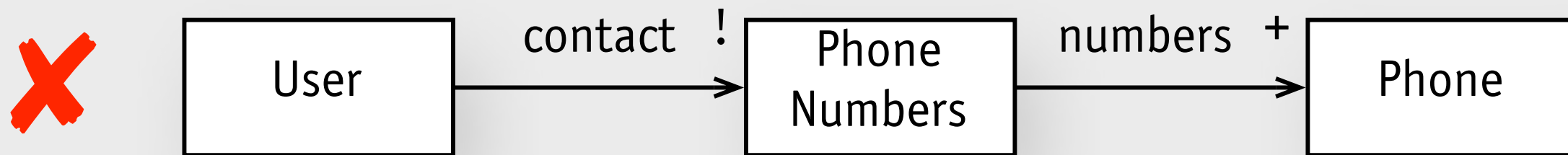friends

User

User

User

Friend

why?
subset is with respect to a context (a user)
without this, data structure won't work

# another example of bad subset

# collections aren't domain objects

✗  | User | —contact !→ | Phone Numbers | —numbers +→ | Phone |

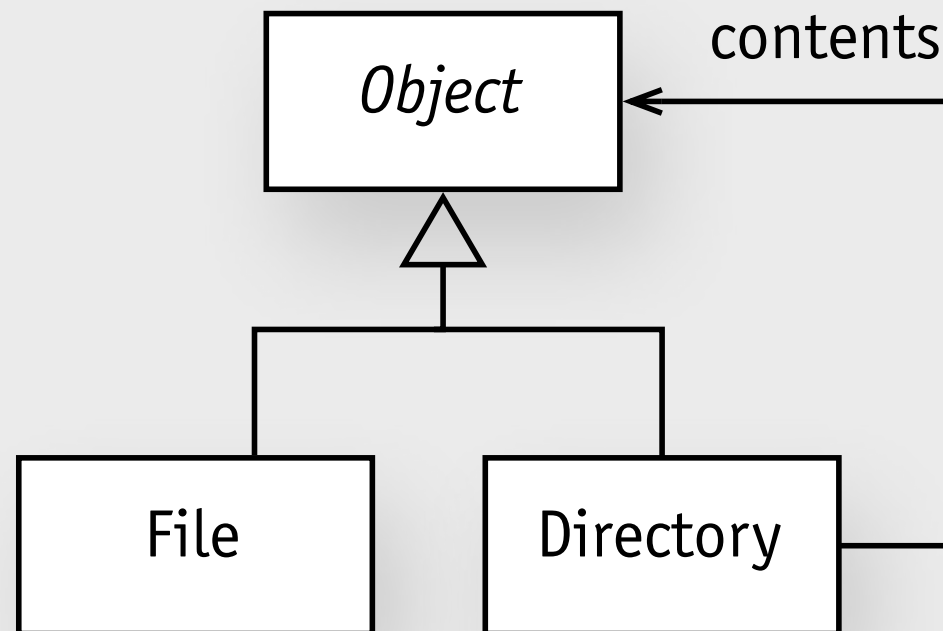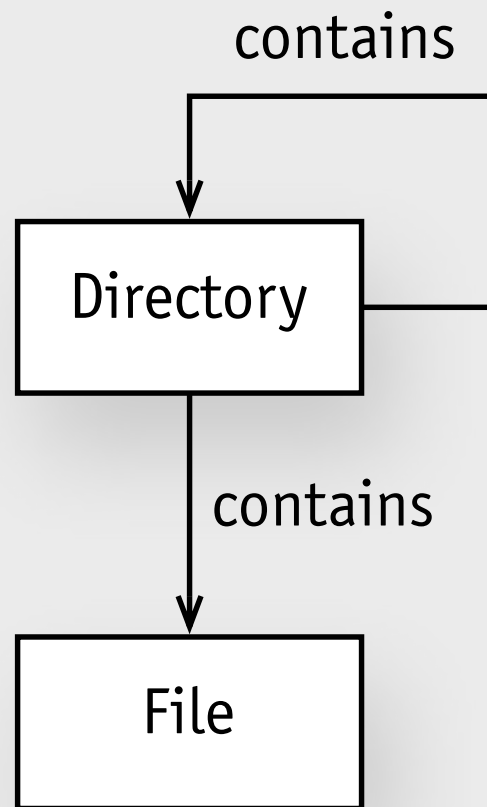✔  | User | —phoneNos +→ | Phone |

✔  | User | —phoneNos +→ | Phone |
        —preferred !→
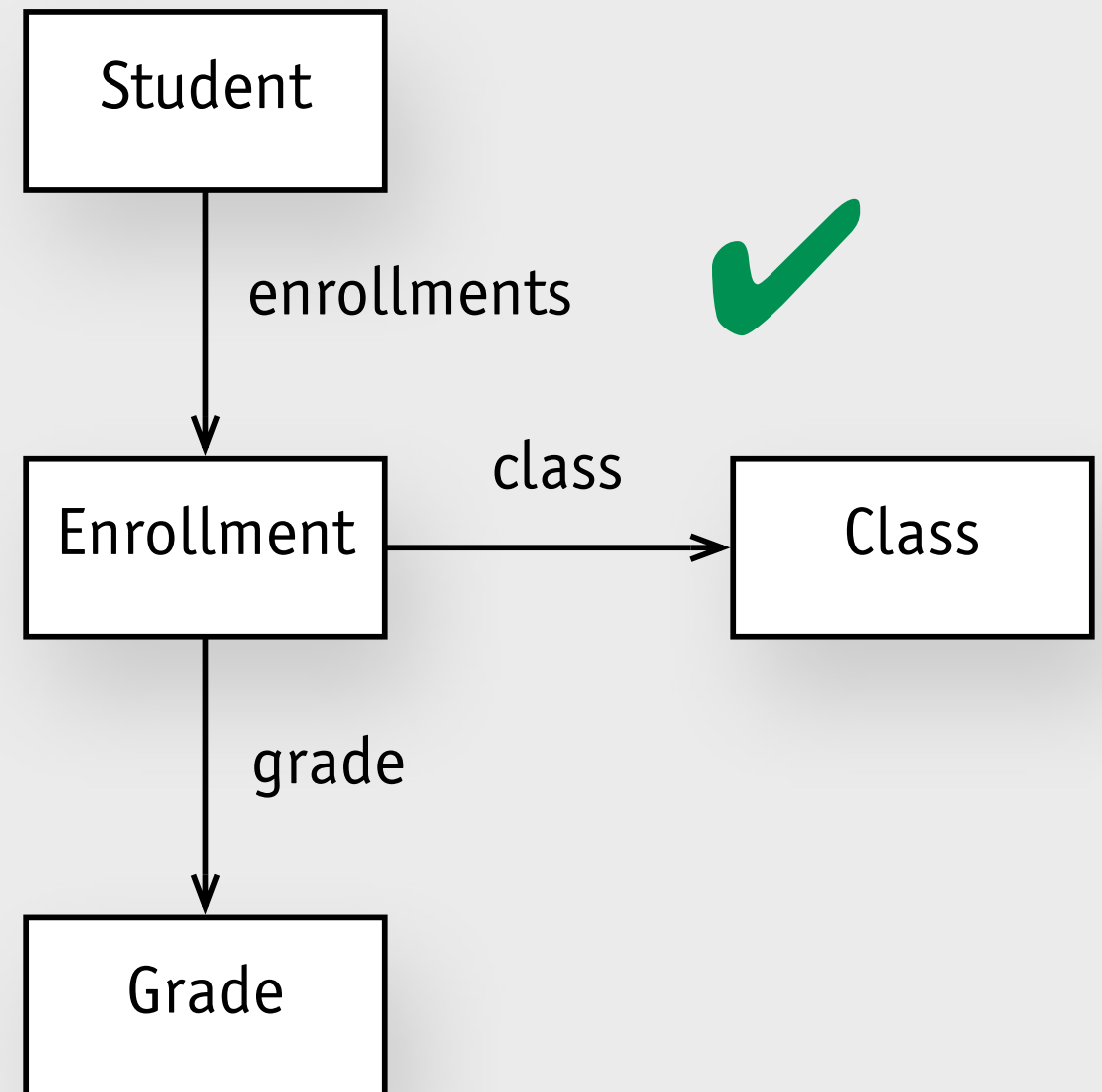
why?
collection objects are implementation details
unless they have properties beyond their contents

# don't split a relation

❌

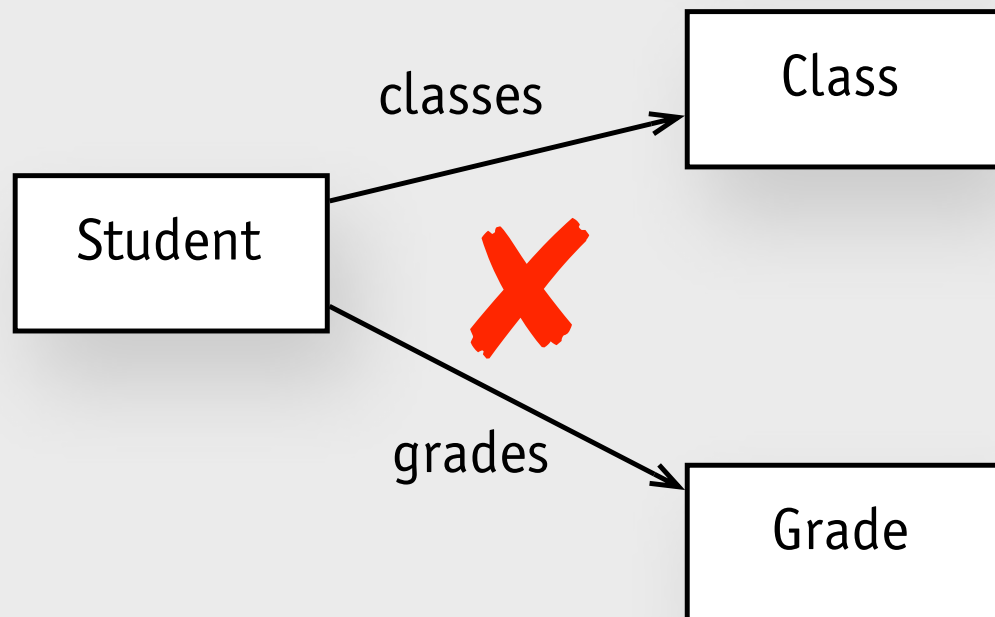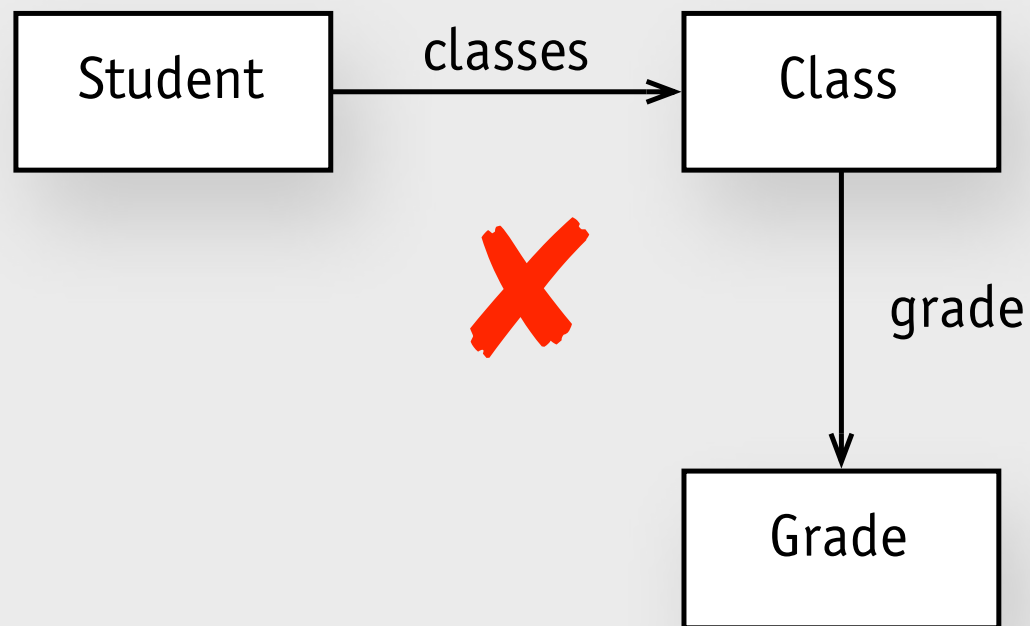contains

Directory

contains

File

✔

Object ←— contents

File    Directory

why?
splitting obscures generalization
leads to duplication in code

# watch out for 3-way relations

Student --classes--> Class

Class --grade--> Grade

❌

Student --classes--> Class

Student --grades--> Grade

❌

Student --enrollments--> Enrollment

Enrollment --class--> Class

Enrollment --grade--> Grade

✔️

why?
student-class-grade is a 3-way relation
need a "tuple" type such as Enrollment

# another example of 3-way relation

# model data that matters

don't ask 'what do I know about students?'
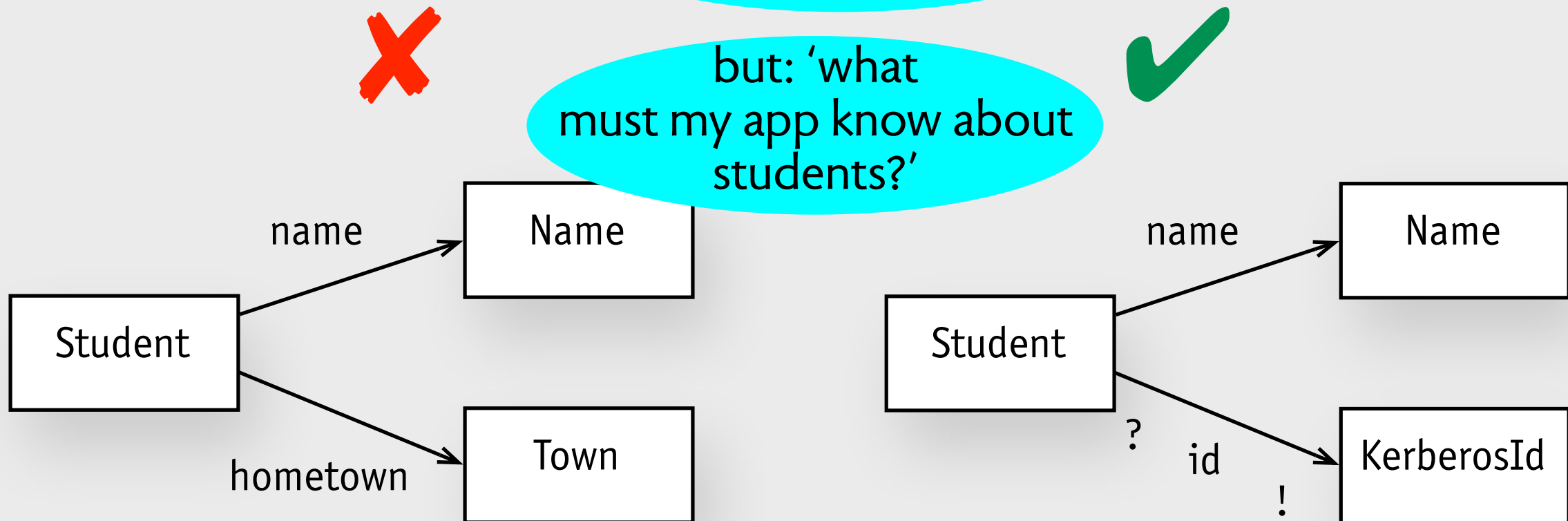
but: 'what must my app know about students?'

❌

✔

**Student** --name--> **Name**

**Student** --hometown--> **Town**

**Student** --name--> **Name**

**Student** --id--> **KerberosId**

?

!

why?
no point modeling the easy stuff
focus on the hard parts
in this case how a student is identified matters
home town probably does not (except maybe for MIT Giving)